

ANALYSIS AND PRACTICAL USE OF FLEXIBLE BICGSTAB

JIE CHEN*, LOIS CURFMAN MCINNES*, AND HONG ZHANG*

Abstract. A flexible version of the BiCGStab algorithm for solving a linear system of equations is analyzed. We show that under variable preconditioning, the perturbation to the outer residual norm is of the same order as that to the application of the preconditioner. Hence, in order to maintain a similar convergence behavior to BiCGStab while reducing the preconditioning cost, the flexible version can be used with a moderate tolerance in the preconditioning Krylov solves. We explored the use of flexible BiCGStab in a large-scale reacting flow application, PFLOTRAN, and showed that the use of a variable multigrid preconditioner significantly accelerates the simulation time on extreme-scale computers using $O(10^4)$ – $O(10^5)$ processor cores.

Key words. BiCGStab, variable preconditioning

AMS subject classifications. 65F10

1. Introduction. Flexible iterative methods [8, 13, 14, 21, 23, 33, 37] for solving a linear system of equations refer to preconditioned Krylov methods where the preconditioner may vary across iterations. The flexible preconditioning strategy is also known under various terms such as *variable* or *inexact preconditioning*. A representative scenario is that the preconditioning requires a linear solve with a second iterative method, in which case “inner iterations” refer to preconditioning and “outer iterations” refer to the flexible Krylov method itself. Flexible methods are an important class of methods that offer several advantages over the use of a fixed preconditioner, one of which is the flexibility to balance the accuracy of the preconditioning solves and the speed of convergence of the outer Krylov iterations in order to reduce the total computational cost. Furthermore, in large-scale applications, the changing landscape of both scientific needs (complex physical models and couplings) and emerging extreme-scale computing systems give rise to practical preconditioners that are hierarchical or nested [5, 6, 9, 15, 18, 26]. Many of these emerging preconditioners benefit from inexact inner solves and thus require the use of outer flexible Krylov methods.

Among many proposed flexible iterative methods, flexible GMRES (abbreviated as FGMRES [23]) is the most frequently used in practice. Its wide use is probably linked to the robustness that results from the long-term recurrence and global orthogonality. Compared with standard GMRES [25], even though the traditional notion of a Krylov subspace is lost, FGMRES computes an orthonormal basis of a subspace within which an optimal residual is sought. Hence, FGMRES still enjoys the residual norm minimization property, and it often shows a satisfactory convergence behavior. On the other hand, in other flexible iterative methods with short-term recurrences, such as inexact PCG [14], flexible QMR [33], and flexible BiCG [37], global (bi)orthogonality is lost, and the convergence behavior is often unpredictable unless the inner solves are sufficiently accurate so that orthogonality is nearly preserved. An idea to improve the robustness is to explicitly perform the orthogonalization as proposed for a variant of the flexible CG algorithm [21]. On the other hand, several analyses of flexible methods, using a larger Krylov subspace that includes the Arnoldi vectors, indicate that the convergence behavior can be maintained with respect to the

*Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439. Emails: (jichen, mcinnes, hzhang)@mcs.anl.gov

fixed preconditioning case if the perturbation to the preconditioner grows inversely with the current residual norm [10, 27, 28].

BiCGStab [34] is a widely used Krylov method. In many applications, BiCGStab outperforms GMRES in terms of both solution time and memory usage, and it has become the de facto method of choice for practitioners. Although BiCGStab is akin to BiCG [11], which generates two sets of biorthogonal residual vectors that naturally form two associated Krylov subspaces, the convergence behavior of BiCGStab is harder to describe because the residual vectors alone do not span the Krylov subspace that contains them. BiCGStab can be seen as redefining the residual polynomial of BiCG by squaring the degree with a smoothing effect. BiCGStab can also be understood as a member of a family of induced dimension reduction (IDR) methods whereby the generated residuals belong to a nested sequence of shrinking subspaces [29, 30, 31], which is analogous to other Krylov methods where there is a similar subspace defined (for example, \mathcal{K}^\perp for CG and $(A\mathcal{K})^\perp$ for GMRES, where \mathcal{K} is the current Krylov subspace). The orthogonal complement notation for the latter methods is consistent with Saad’s viewpoint of projection methods [24]; however, BiCGStab does not belong to this class.

We study the flexible BiCGStab algorithm (FBiCGStab) in this paper. FBiCGStab was initially proposed in [37], and it was recently cast under the framework of flexible variants of IDR methods [35]. However, little is known about its convergence properties. The goal of this paper is to analyze the behavior of the method and to provide guidance on its practical use. We do not study the convergence guarantee of the method, but we argue that the convergence behavior is close to that of the fixed preconditioning case if the perturbation to the preconditioner is not too large. Thus, an appropriate perturbation is key to the performance: if too large, the residual behavior of the outer iterations is unpredictable; if too small, the inner solves may be time consuming to complete. Because of the loose connection of BiCGStab with the associated Krylov subspace, this analysis is different from that for other flexible Krylov methods (see [10, 27, 28]). Rather, the arguments are made on orthogonality and norm minimization properties that guarantee bounds of perturbations on the iterates. Interestingly, this result can lead to a conclusion similar to that in [10, 27, 28], that is, the convergence behavior of the flexible method can be maintained by relaxing the accuracy requirement of the preconditioning solves as the outer residual norm decreases, although this result might be less practically useful (see Section 2.2).

What motivates our study of FBiCGStab is the scaling difficulty of the simulation of reacting flows in a geological application PFLOTRAN [17, 19] on extreme-scale parallel computers. A well-known bottleneck in Krylov methods for solving large-scale linear systems on parallel computers is the inner-product calculations, because they require global synchronizations. The scaling of the solver starts to deteriorate when the number of processor cores increases beyond $O(10^5)$ [2]. Among several strategies, such as deferring or pipelining inner product calculations, overlapping communications with computations, and block orthogonalizations [7, 12, 20, 32, 36, 38], the use of a hierarchical preconditioner that reduces the length of vectors in inner-product calculations has been demonstrated to be effective for inner-outer GMRES iterations [18]. This discovery prompted us to generalize the idea to flexible BiCGStab iterations, because BiCGStab has historically been the preferred solver for PFLOTRAN by domain scientists. Based on the perturbation analysis, we use multigrid preconditioners and see remarkable improvement in the simulation time using $O(10^4)$ – $O(10^5)$ processor cores.

The rest of the paper is organized as follows. With a brief derivation of FBiCGStab, Section 2 analyzes the behavior of the residual norm under flexible preconditioning and gives examples to illustrate the interpretation of the results. Then, several numerical examples are shown in Section 3 to demonstrate the need for a stopping criterion with a moderate tolerance in the preconditioning inner solves. Section 4 presents the successful use of FBiCGStab with multigrid preconditioners in PFLOTRAN. Concluding remarks are given in Section 5.

2. Algorithm and analysis. The following is the standard unpreconditioned BiCGStab algorithm for solving a linear system

$$Ax = b,$$

using x_0 as the initial vector [24]:

```

1:  $r_0 = b - Ax_0$ ;  $\bar{r}_0$  arbitrary but  $(\bar{r}_0, r_0) \neq 0$ 
2:  $p_0 = r_0$ 
3: for  $j = 0, 1, \dots$  until convergence do
4:    $\alpha_j = (r_j, \bar{r}_0) / (Ap_j, \bar{r}_0)$ 
5:    $s_j = r_j - \alpha_j Ap_j$ 
6:    $\omega_j = (As_j, s_j) / (As_j, As_j)$ 
7:    $x_{j+1} = x_j + \alpha_j p_j + \omega_j s_j$ 
8:    $r_{j+1} = s_j - \omega_j As_j$ 
9:    $\beta_j = (r_{j+1}, \bar{r}_0) / (r_j, \bar{r}_0) \cdot \alpha_j / \omega_j$ 
10:   $p_{j+1} = r_{j+1} + \beta_j (p_j - \omega_j Ap_j)$ 
11: end for
```

The coefficient iterates α_j and β_j are derived based on their counterparts in BiCG for updating the residual vectors and the search direction vectors. One can show that α_j makes $s_j \perp \bar{r}_0$ and β_j makes $p_{j+1} \perp \bar{r}_0$ for all j . Furthermore, ω_j is defined to minimize the 2-norm of the residual vector r_{j+1} given s_j and As_j .

When the algorithm is used with a preconditioner $M \approx A$, the right preconditioning consists of solving the system

$$AM^{-1}y = b, \quad y = Mx.$$

One way to derive the preconditioned iteration is, in the above unpreconditioned version, to replace the symbol A by AM^{-1} and x_j by y_j , and then substitute y_j back by Mx_j . This introduces two auxiliary vectors $\tilde{p}_j = M^{-1}p_j$ and $\tilde{s}_j = M^{-1}s_j$, which are the only computations that require the preconditioner. We summarize this preconditioned version in Algorithm 1. It is the same as the one presented in [37].

One would also like to consider left preconditioning, where M^{-1} is applied to the left of the system $Ax = b$. After a change of variables, the left preconditioned version is almost the same as Algorithm 1, except that the two inner products in line 8 are changed to the ones using $(MM^T)^{-1}$ -norm. In this case, ω_j minimizes the $(MM^T)^{-1}$ -norm of r_{j+1} . Compared with right preconditioning, left preconditioning incurs two more applications of the preconditioner in each iteration, which increases the computational cost. Furthermore, similar to other Krylov methods, a major hurdle for developing variable preconditioners for left preconditioning is the disconnection between the preconditioned residuals and the actual residuals. Therefore, we do not consider left preconditioning in this paper.

Algorithm 1 Right preconditioned BiCGStab / Flexible BiCGStab

```

1:  $r_0 = b - Ax_0$ ;  $\bar{r}_0$  arbitrary
2:  $p_0 = r_0$ 
3: for  $j = 0, 1, \dots$  until convergence do
4:    $\tilde{p}_j = M^{-1}p_j$ 
5:    $\alpha_j = (r_j, \bar{r}_0)/(A\tilde{p}_j, \bar{r}_0)$ 
6:    $s_j = r_j - \alpha_j A\tilde{p}_j$ 
7:    $\tilde{s}_j = M^{-1}s_j$ 
8:    $\omega_j = (A\tilde{s}_j, s_j)/(A\tilde{s}_j, A\tilde{s}_j)$ 
9:    $x_{j+1} = x_j + \alpha_j \tilde{p}_j + \omega_j \tilde{s}_j$ 
10:   $r_{j+1} = s_j - \omega_j A\tilde{s}_j$ 
11:   $\beta_j = (r_{j+1}, \bar{r}_0)/(r_j, \bar{r}_0) \cdot \alpha_j / \omega_j$ 
12:   $p_{j+1} = r_{j+1} + \beta_j(p_j - \omega_j A\tilde{p}_j)$ 
13: end for

```

An iterative method can be used to compute \tilde{p}_j in line 4 and \tilde{s}_j in line 7 of Algorithm 1, but the iterations may not run to full accuracy. In this case, Algorithm 1 becomes the flexible version of BiCGStab. The computed iterates \tilde{p}_j and \tilde{s}_j under inexact preconditioning will carry on their error to subsequent iterations. To gauge the amplification of error, we are interested in the situation that the relative residual of the inner solves with M is bounded by a small tolerance ϵ . That is, if we use an underline to denote the actual iterates with errors, we assume that

$$\|\underline{p}_j - M\underline{\tilde{p}}_j\| \leq \epsilon \|\underline{p}_j\| \quad \text{and} \quad \|\underline{s}_j - M\underline{\tilde{s}}_j\| \leq \epsilon \|\underline{s}_j\|. \quad (2.1)$$

The following subsection characterizes the relative difference between r_{j+1} and \underline{r}_{j+1} under condition (2.1).

2.1. Analysis. The analysis is based on the fact that the coefficient iterates α_j , β_j , and ω_j are computed such that the inaccuracy incurred in the preconditioning solves is not “adversely” accumulated to affect outer iterations. For this, we need the following observations. They are trivially correct in the fixed preconditioned case, but they are also true when a flexible preconditioner is used. The proof is simple and thus omitted.

PROPOSITION 2.1. *The iterates in Algorithm 1 have the following properties:*

- (i) *The vector r_{j+1} is the residual, that is, $r_{j+1} = b - Ax_{j+1}$.*
- (ii) *Consider that s_j is a function of α_j ; then the definition of α_j in line 5 makes $s_j \perp \bar{r}_0$.*
- (iii) *Consider that p_{j+1} is a function of β_j ; then the definition of β_j in line 11 makes $p_{j+1} \perp \bar{r}_0$.*
- (iv) *Consider that r_{j+1} is a function of ω_j ; then the definition of ω_j in line 8 minimizes $\|r_{j+1}\|_2$.*

In light of these observations, we will bound the perturbation to vectors in the form $x - \alpha y$, where the scalar α is used to satisfy some orthogonality or minimization property when x and y are perturbed. We will make heavy use of angles spanned by two \mathbb{R}^n vectors. For this, we use $\angle(x, y)$ to denote the *acute* angle between the two vectors, that is, $|(x, y)| = \|x\|\|y\|\cos\angle(x, y)$. Hence, $\cos\angle(x, y)$ is always nonnegative. The following lemma establishes the basic fact, generalizing from the intuitive \mathbb{R}^3 case, that for three vectors spanning three angles, one of the angles is bounded by the

sum and the difference of the other two. Then, two lemmas follow, stating that the perturbation to $x - \alpha y$ is in the same order as that to x and y .

LEMMA 2.2. *If all pairwise angles among three vectors x , y , and z are acute, then*

$$|\angle(x, y) - \angle(y, z)| \leq \angle(x, z) \leq \angle(x, y) + \angle(y, z).$$

Proof. Without loss of generality, we assume that $\|x\| = \|y\| = \|z\| = 1$. It suffices to prove, for the case $\angle(x, y) + \angle(y, z)$ is acute, that

$$\begin{aligned} \cos \angle(x, y) \cos \angle(y, z) - \sin \angle(x, y) \sin \angle(y, z) &\leq \cos \angle(x, z) \leq \\ &\cos \angle(x, y) \cos \angle(y, z) + \sin \angle(x, y) \sin \angle(y, z), \end{aligned}$$

which is equivalent to

$$|(x, z) - (x, y)(y, z)| \leq \sqrt{1 - (x, y)^2} \sqrt{1 - (y, z)^2}. \quad (2.2)$$

Note that

$$|(x, z) - (x, y)(y, z)| = |x^T(I - yy^T)z|$$

and that

$$\sqrt{1 - (x, y)^2} \sqrt{1 - (y, z)^2} = \sqrt{x^T(I - yy^T)x} \sqrt{z^T(I - yy^T)z}.$$

Thus, Cauchy's inequality for vector semi-inner products with respect to the symmetric positive semi-definite matrix $I - yy^T$ proves (2.2). \square

LEMMA 2.3. *Given a vector r , let $z = x - \alpha y$ and $\underline{z} = \underline{x} - \underline{\alpha} y$, where $\alpha = (x, r)/(y, r)$ and $\underline{\alpha} = (\underline{x}, r)/(\underline{y}, r)$, and let $\gamma = \text{sgn}[(x, y)(x, r)(y, r)]$. If there exist ϵ_x, ϵ_y such that $\epsilon_x < \cos \angle(x, r)$, $\epsilon_y < \cos \angle(y, r)$ and that*

$$\|x - \underline{x}\| \leq \epsilon_x \|x\| \quad \text{and} \quad \|y - \underline{y}\| \leq \epsilon_y \|y\|,$$

then $\|z - \underline{z}\| \leq \epsilon_z \|z\|$ with

$$\begin{aligned} \epsilon_z &= \frac{\epsilon_x + \frac{\cos \angle(x, r)}{\cos \angle(y, r)} [2\epsilon_y + (1 + \epsilon_x)(1 + B)(1 + D) - 1]}{\sqrt{1 + \frac{\cos^2 \angle(x, r)}{\cos^2 \angle(y, r)} - 2\gamma \cos \angle(x, y) \frac{\cos \angle(x, r)}{\cos \angle(y, r)}}}, \quad (2.3) \\ B &= \epsilon_x \left(\frac{\epsilon_x}{2\sqrt{1 - \epsilon_x^2}} + \tan \angle(x, r) \right), \quad D = \frac{2\epsilon_y \sqrt{1 - \epsilon_y^2} \tan \angle(y, r) + \epsilon_y^2}{-2\epsilon_y \sqrt{1 - \epsilon_y^2} \tan \angle(y, r) + 2(1 - \epsilon_y^2)}. \end{aligned}$$

In other words, denoting by $\epsilon = \max\{\epsilon_x, \epsilon_y\}$, we have $\epsilon_z \leq C\epsilon$ where the prefactor C has a finite limit

$$\lim_{\epsilon \rightarrow 0} C = \frac{\cos \angle(y, r) + \cos \angle(x, r) [3 + \tan \angle(x, r) + \tan \angle(y, r)]}{\sqrt{\cos^2 \angle(y, r) + \cos^2 \angle(x, r) - 2\gamma \cos \angle(x, y) \cos \angle(x, r) \cos \angle(y, r)}}.$$

Proof. First we have

$$\begin{aligned} \|z - \underline{z}\| &\leq \|x - \underline{x}\| + \|\alpha y - \underline{\alpha y}\| \\ &\leq \|x - \underline{x}\| + \|\alpha(y - \underline{y})\| + \|(\alpha - \underline{\alpha})\underline{y}\| \\ &\leq \epsilon_x \|x\| + \epsilon_y \|\alpha y\| + \|(\alpha - \underline{\alpha})\underline{y}\|. \end{aligned}$$

Based on Lemma 2.2, when $\epsilon_x < \cos \angle(x, r)$, the angle between x and r and that between \underline{x} and r will be acute, or obtuse, at the same time. This means that the inner products (x, r) and (\underline{x}, r) have the same sign. Similarly it holds for (y, r) and (\underline{y}, r) . Thus, α and $\underline{\alpha}$ have the same sign, and hence

$$\frac{\|(\alpha - \underline{\alpha})\underline{y}\|}{\|\alpha y\|} = \left| \frac{\|\underline{y}\|}{\|y\|} - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right|.$$

Using the fact that $|\|\underline{y}\|/\|y\| - 1| \leq \epsilon_y$, we get

$$\|(\alpha - \underline{\alpha})\underline{y}\| \leq \|\alpha y\| \left(\epsilon_y + \left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| \right). \quad (2.4)$$

With $\|\alpha y\| = \|x\| \cos \angle(x, r) / \cos \angle(y, r)$ and

$$\|z\|^2 = \|x\|^2 \left(1 + \frac{\cos^2 \angle(x, r)}{\cos^2 \angle(y, r)} - 2\gamma \cos \angle(x, y) \frac{\cos \angle(x, r)}{\cos \angle(y, r)} \right),$$

we thus obtain

$$\begin{aligned} \|z - \underline{z}\| &\leq \epsilon_x \|x\| + \|\alpha y\| \left(2\epsilon_y + \left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| \right) \\ &= \|z\| \frac{\epsilon_x + \frac{\cos \angle(x, r)}{\cos \angle(y, r)} \left(2\epsilon_y + \left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| \right)}{\sqrt{1 + \frac{\cos^2 \angle(x, r)}{\cos^2 \angle(y, r)} - 2\gamma \cos \angle(x, y) \frac{\cos \angle(x, r)}{\cos \angle(y, r)}}}. \end{aligned} \quad (2.5)$$

Therefore, we proceed to bound

$$\left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| = \left| 1 - \frac{\|\underline{x}\| \cos \angle(\underline{x}, r) \cos \angle(y, r)}{\|x\| \cos \angle(x, r) \cos \angle(\underline{y}, r)} \right|.$$

The strategy of bounding this term is to find $A, B, D > 0$ such that

$$\left| 1 - \frac{\|\underline{x}\|}{\|x\|} \right| \leq A, \quad \left| 1 - \frac{\cos \angle(\underline{x}, r)}{\cos \angle(x, r)} \right| \leq B, \quad \left| 1 - \frac{\cos \angle(y, r)}{\cos \angle(\underline{y}, r)} \right| \leq D;$$

then,

$$\left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| \leq (1 + A)(1 + B)(1 + D) - 1. \quad (2.6)$$

Clearly, we can let $A = \epsilon_x$.

To simplify notation, let $\angle(x, r) = \beta$ and $\angle(\underline{x}, r) = \beta + \delta$ for some δ . Then,

$$\begin{aligned} \left| 1 - \frac{\cos(\beta + \delta)}{\cos \beta} \right| &= |(1 - \cos \delta) + \sin \delta \tan \beta| \\ &= \left| (\sin \delta) \left(\tan \frac{\delta}{2} + \tan \beta \right) \right| \leq |\sin \delta| \left(\frac{1}{2} |\tan \delta| + \tan \beta \right). \end{aligned} \quad (2.7)$$

Because $|\sin \delta| \leq \epsilon_x$,

$$\left| 1 - \frac{\cos(\beta + \delta)}{\cos \beta} \right| \leq \epsilon_x \left(\frac{\epsilon_x}{2\sqrt{1 - \epsilon_x^2}} + \tan \angle(x, r) \right) =: B. \quad (2.8)$$

Similarly, let $\angle(y, r) = \eta$ and $\angle(\underline{y}, r) = \eta + \tau$ for some τ . Note that based on Lemma 2.2, $|\tau| \leq \pi/2 - \eta$. Using

$$\left| \tan \frac{\tau}{2} \right| \leq \frac{1}{2} |\tan \tau| \leq \frac{\epsilon_y}{2\sqrt{1 - \epsilon_y^2}},$$

we have

$$\left| 1 - \frac{\cos \eta}{\cos(\eta + \tau)} \right| = \left| \frac{\tan \eta + \tan(\tau/2)}{\tan \eta - \cot \tau} \right| \leq \frac{\tan \angle(y, r) + \frac{\epsilon_y}{2\sqrt{1 - \epsilon_y^2}}}{-\tan \angle(y, r) + \frac{\sqrt{1 - \epsilon_y^2}}{\epsilon_y}} =: D. \quad (2.9)$$

With (2.8) and (2.9), we establish (2.6). Then, together with (2.5), we have proved (2.3). \square

LEMMA 2.4. *Let $z = x - \alpha y$ and $\underline{z} = \underline{x} - \underline{\alpha} y$, where $\alpha = (x, y)/\|y\|^2$ and $\underline{\alpha} = (\underline{x}, \underline{y})/\|\underline{y}\|^2$. If (x, y) and $(\underline{x}, \underline{y})$ have the same sign and there exist ϵ_x, ϵ_y such that $\epsilon_x + \epsilon_y < 1$ and that*

$$\|x - \underline{x}\| \leq \epsilon_x \|x\| \quad \text{and} \quad \|y - \underline{y}\| \leq \epsilon_y \|y\|,$$

then $|\alpha - \underline{\alpha}| \leq \epsilon_\alpha |\alpha|$ and $\|z - \underline{z}\| \leq \epsilon_z \|z\|$ with

$$\epsilon_\alpha = (\epsilon_x + \epsilon_y)(1 + \epsilon_y) \left\{ 1 + (1 + \epsilon_x) \left[\frac{\epsilon_x + \epsilon_y}{2\sqrt{1 - (\epsilon_x + \epsilon_y)^2}} + \tan \angle(x, y) \right] \right\} \quad (2.10)$$

and

$$\epsilon_z = (\epsilon_x + \epsilon_y)(1 + \epsilon_x) + \frac{\epsilon_x}{\sin \angle(x, y)} + \frac{\epsilon_y}{\tan \angle(x, y)} + \frac{\epsilon_x + \epsilon_y}{\tan \angle(x, y)} \left[1 + \frac{(1 + \epsilon_x)(\epsilon_x + \epsilon_y)}{2\sqrt{1 - (\epsilon_x + \epsilon_y)^2}} \right]. \quad (2.11)$$

In other words, denoting by $\epsilon = \max\{\epsilon_x, \epsilon_y\}$, we have $\epsilon_\alpha \leq C_\alpha \epsilon$ and $\epsilon_z \leq C_z \epsilon$, where the prefactors C_α and C_z have finite limits

$$\lim_{\epsilon \rightarrow 0} C_\alpha = 2 + 2 \tan \angle(x, y), \quad \lim_{\epsilon \rightarrow 0} C_z = \left(2 + \frac{1 + 3 \cos \angle(x, y)}{\sin \angle(x, y)} \right).$$

Proof. Using a similar argument as in the proof of the preceding lemma, we reach

$$\begin{aligned} \|z - \underline{z}\| &\leq \epsilon_x \|x\| + \epsilon_y \|\alpha y\| + \|(\alpha - \underline{\alpha}) y\| \\ &\leq \|z\| \left[\frac{\epsilon_x}{\sin \angle(x, y)} + \cot \angle(x, y) \left(2\epsilon_y + \left| 1 - \frac{\|\alpha y\|}{\|\alpha y\|} \right| \right) \right], \end{aligned} \quad (2.12)$$

by noting that $\|\alpha y\| = \|x\| \cos \angle(x, y)$, that $\|z\| = \|x\| \sin \angle(x, y)$, and that α and $\underline{\alpha}$ have the same sign by the condition of this lemma. Furthermore, with (2.4), which clearly also holds here, we have

$$\frac{|\alpha - \underline{\alpha}|}{|\alpha|} \leq \frac{\|y\|}{\|\underline{y}\|} \left(\epsilon_y + \left| 1 - \frac{\|\alpha y\|}{\|\alpha y\|} \right| \right) \leq (1 + \epsilon_y) \left(\epsilon_y + \left| 1 - \frac{\|\alpha y\|}{\|\alpha y\|} \right| \right). \quad (2.13)$$

Therefore, we proceed to bound

$$\left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| = \left| 1 - \frac{\|\underline{x}\| \cos \angle(\underline{x}, y)}{\|x\| \cos \angle(x, y)} \right|.$$

The strategy of bounding this term is to find $A, B > 0$ such that

$$\left| 1 - \frac{\|\underline{x}\|}{\|x\|} \right| \leq A, \quad \left| 1 - \frac{\cos \angle(\underline{x}, y)}{\cos \angle(x, y)} \right| \leq B;$$

then,

$$\left| 1 - \frac{\|\underline{\alpha y}\|}{\|\alpha y\|} \right| \leq (1 + A)(1 + B) - 1. \quad (2.14)$$

Clearly, we can let $A = \epsilon_x$.

To simplify notation, let $\angle(x, y) = \beta$ and $\angle(\underline{x}, y) = \beta + \delta$ for some δ . Then, the same as (2.7), we have

$$\left| 1 - \frac{\cos(\beta + \delta)}{\cos \beta} \right| \leq |\sin \delta| \left(\frac{1}{2} |\tan \delta| + \tan \beta \right).$$

Let θ_x be the angle between x and \underline{x} , and similarly for θ_y . Because $\|x - \underline{x}\| \leq \epsilon_x \|x\|$ with $\epsilon_x < 1$, θ_x is acute. Similarly, so is θ_y . Note that $\|x\| \sin \theta_x \leq \|x - \underline{x}\| \leq \epsilon_x \|x\|$; therefore, $\sin \theta_x \leq \epsilon_x$, and similarly $\sin \theta_y \leq \epsilon_y$. Then, $\sin \theta_x \leq \epsilon_x < \sqrt{1 - \epsilon_y^2} \leq \cos \theta_y$, which indicates that $\theta_x + \theta_y$ is also acute. Thus, the fact that $|\delta| \leq \theta_x + \theta_y$ leads to $\sin |\delta| \leq \sin \theta_x + \sin \theta_y \leq \epsilon_x + \epsilon_y$. Therefore

$$\left| 1 - \frac{\cos(\beta + \delta)}{\cos \beta} \right| \leq (\epsilon_x + \epsilon_y) \left[\frac{\epsilon_x + \epsilon_y}{2\sqrt{1 - (\epsilon_x + \epsilon_y)^2}} + \tan \angle(x, y) \right] =: B. \quad (2.15)$$

With (2.15), we establish (2.14). Then, together with (2.13) and (2.12), we have proved (2.10) and (2.11). \square

Using the above two lemmas, we have the following result. It states that the relative perturbation to the outer residual norm is in the same order as the relative residual norm in the inner solves.

THEOREM 2.5. *If Algorithm 1 (with flexible preconditioning) is run with a finite number of iterations where neither breakdown nor stagnation occurs, then under condition (2.1), for all j ,*

$$\frac{\|r_j - \underline{r}_j\|}{\|r_j\|} = O(\epsilon). \quad (2.16)$$

Proof. To facilitate presentation, we define $\text{err}(a) := \|a - \underline{a}\|/\|a\|$ for any vector or scalar $a \neq 0$. We first observe that

$$\begin{aligned} \frac{\|A\tilde{p}_j - A\underline{\tilde{p}}_j\|}{\|A\tilde{p}_j\|} &= \frac{\|AM^{-1}(p_j - M\underline{\tilde{p}}_j)\|}{\|AM^{-1}p_j\|} \leq \kappa \frac{\|p_j - M\underline{\tilde{p}}_j\|}{\|p_j\|} \\ &\leq \kappa \left(\frac{\|p_j - \underline{p}_j\|}{\|p_j\|} + \frac{\|\underline{p}_j - M\underline{\tilde{p}}_j\|}{\|p_j\|} \right) \leq \kappa \left(\frac{\|p_j - \underline{p}_j\|}{\|p_j\|} + \epsilon \frac{\|\underline{p}_j\|}{\|p_j\|} \right), \end{aligned}$$

where κ denotes the condition number of AM^{-1} . Since the big-O notation is used for sufficiently small ϵ , the above observation means that if $\text{err}(p_j) = O(\epsilon)$, then $\text{err}(A\tilde{p}_j) = O(\epsilon)$. Similarly, if $\text{err}(s_j) = O(\epsilon)$, then $\text{err}(A\tilde{s}_j) = O(\epsilon)$.

We now show the theorem by induction on $\text{err}(r_j)$ and $\text{err}(p_j)$. The conditions in the theorem (no breakdown and stagnation) are used to ensure the applicability of the lemmas. At $j = 0$, r_0 and p_0 are unchanged under variable preconditioning. If $\text{err}(r_j) = O(\epsilon)$ and $\text{err}(p_j) = O(\epsilon)$, then because $\text{err}(A\tilde{p}_j) = O(\epsilon)$, we have $\text{err}(s_j) = O(\epsilon)$ by Lemma 2.3. Consequently, $\text{err}(A\tilde{s}_j) = O(\epsilon)$, and thus $\text{err}(r_{j+1}) = O(\epsilon)$ by Lemma 2.4.

Now consider $p_{j+1} = r_{j+1} + \beta_j z_j$ where $z_j = p_j - \omega_j A\tilde{p}_j$. Both $\text{err}(p_j)$ and $\text{err}(A\tilde{p}_j)$ are $O(\epsilon)$. On the other hand, $\text{err}(\omega_j)$ is also $O(\epsilon)$ according to Lemma 2.4, because we have shown that both $\text{err}(s_j)$ and $\text{err}(A\tilde{s}_j)$ are $O(\epsilon)$. Hence, $\text{err}(z_j) = O(\epsilon)$. Then by invoking Lemma 2.3 again we have $\text{err}(p_{j+1}) = O(\epsilon)$, which completes the induction. \square

Theorem 2.5 considers the perturbation of the residual in the relative sense. As a corollary, a result for the absolute perturbation is given next. Instead of a fixed tolerance ϵ for all the inner solves, we allow the tolerance, denoted by ϵ_j , to vary in each outer iteration j . The result indicates a reciprocal relationship between the residual norm $\|r_j\|$ and ϵ_j .

COROLLARY 2.6. *Let the relative inner tolerance ϵ depend on the outer iterations indexed by j , that is,*

$$\|p_j - M\tilde{p}_j\| \leq \epsilon_j \|p_j\| \quad \text{and} \quad \|s_j - M\tilde{s}_j\| \leq \epsilon_j \|s_j\|.$$

Under the condition of Theorem 2.5, if the residual norm $\|r_j\|$ is monotonically decreasing, then for any δ there exists a constant C such that if

$$\epsilon_j = \frac{C\delta}{\|r_j\|},$$

then $\|r_j - r_j\| \leq \delta$.

Proof. Note that Theorem 2.5 is proved by induction on j . When ϵ_j is monotonically increasing but is still sufficiently small (guaranteed by a finite number of iterations), a stronger conclusion is that there exists a C' that is independent of j such that

$$\text{err}(r_j) \leq C' \epsilon_j, \tag{2.17}$$

because in the right-hand side ϵ_j can always be relaxed later by changing it to ϵ_{j+1} . Rewriting (2.17), we have $\|r_j - r_j\| \leq C' \epsilon_j \|r_j\|$. Therefore, if we let $\epsilon_j = C\delta/\|r_j\|$ by using some C such that $C'C \leq 1$ and that $C\delta/\|r_j\|$ is sufficiently small to trigger the validity of (2.17), we immediately have $\|r_j - r_j\| \leq \delta$. \square

2.2. Interpretation and illustrative examples. Roughly speaking, the above results state that if the perturbation to the preconditioning step is not relatively large, then the convergence history looks similar to that of the case without perturbation. Whereas this gives a qualitative description of the behavior of BiCGStab under variable preconditioning, one must be cautious in interpreting the theoretical results and not overly emphasize their predictive power in a quantitative manner. First and most important, these results are not convergence assertions. Since BiCGStab itself is not guaranteed to converge for general unsymmetric matrices, it does not seem reasonable to expect that a flexible version magically proves the opposite.

Second, the perturbation result is built based on a small- ϵ regime. Clearly, when $\epsilon \rightarrow 0$, such as when ϵ is the machine precision, the term “flexible” gradually loses its meaning since the convergence history will be almost identical to that of the fixed preconditioning case. A tricky question is determining when the behavior of flexible BiCGStab starts to diverge from that of standard BiCGStab, since it depends on when the accumulated prefactor in front of ϵ , as in (2.16), grows to an unacceptable level. It is unclear how to characterize this prefactor, but the prefactor is certainly connected to the bound ϵ_z in Lemma 2.3 and the bounds ϵ_α and ϵ_z in Lemma 2.4. Empirically we see that the perturbation in the range of $O(10^{-4})$ – $O(10^{-2})$ gives reasonably similar convergence behavior to the case of no perturbation ($\epsilon = 0$).

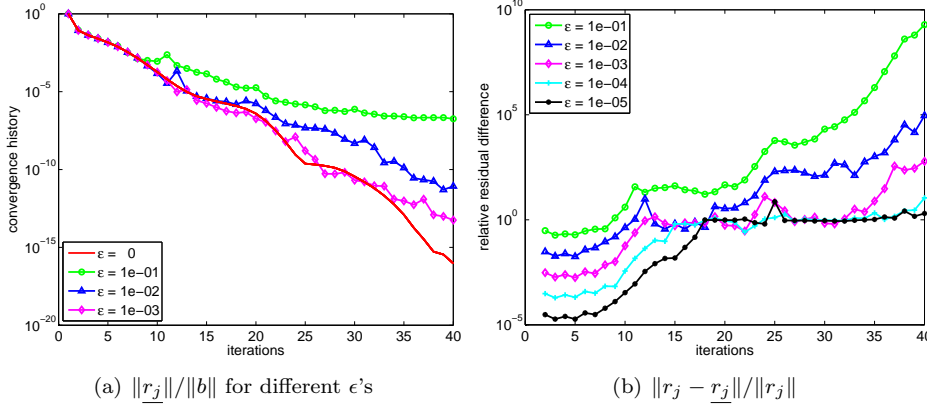


FIG. 2.1. Convergence history and $\|r_j - \underline{r}_j\|/\|r_j\|$ under a fixed level of perturbation ϵ in the preconditioning step. In (a) a few curves are omitted to avoid cluttering.

Figure 2.1 shows the convergence results of solving a model problem that is introduced in more detail in Section 3. Here, the discretization of the problem in 2D yields a matrix of size $65,536 \times 65,536$, and a choice of the parameters $\gamma = 4/h$ and $\beta = -0.2/h^2$ renders the matrix positive definite but unsymmetric. With block Jacobi/ILU(0) preconditioning of block size 64×64 and a zero initial vector, the residuals monotonically decrease to machine precision in 40 iterations (see the red curve without markers in plot (a)).

We artificially perturbed the vectors p_j and s_j (see (2.1)) by a Gaussian random vector normalized to a norm of $\epsilon\|p_j\|$ and $\epsilon\|s_j\|$, respectively, before computing the preconditioned vectors \tilde{p}_j and \tilde{s}_j by block Jacobi/ILU(0). This mimics the situation of an inexact preconditioning solve with relative residual tolerance ϵ . One sees that in plot (a), as ϵ decreases, the convergence history is closer and closer to the reference red curve. We also ran experiments with smaller ϵ 's, but the residual curves were so close to the reference curve that we omit them in the plot to avoid cluttering. In plot (b), we show $\|r_j - \underline{r}_j\|/\|r_j\|$. As predicted by Theorem 2.5, the curves corresponding to a smaller ϵ tend to be positioned below those corresponding to a larger ϵ .

The absolute difference, $\|r_j - \underline{r}_j\|$, on the other hand, can be made bounded as opposed to the increasing trends of Figure 2.1(b). One simple way is, in fact, to do nothing, because $\|r_j\|$ itself decreases. More interesting is that we can increase ϵ across iterations, as Corollary 2.6 implies. In Figure 2.2 we define $\epsilon_j = D/\|r_j\|$ for various different choices of D . The iterations converge to a level dependent on D .

One is tempted to infer from Corollary 2.6 that successively relaxing the pertur-

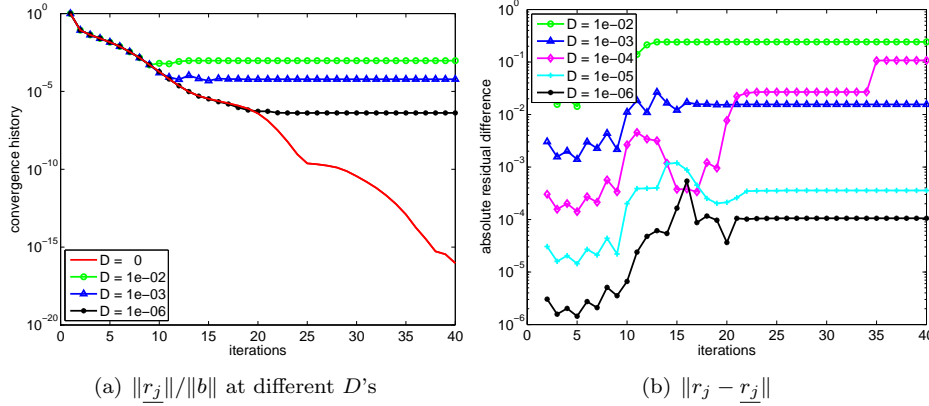


FIG. 2.2. Convergence history and $\|r_j - \underline{r}_j\|/\|r_j\|$ under a dynamically adjusted perturbation $\epsilon_j = D/\|r_j\|$ in the preconditioning step. In (a) a few curves are omitted to avoid cluttering.

bation ϵ_j in this manner can perhaps reduce the preconditioning cost as the iterations converge. However, several caveats render this corollary to be less practically useful than Theorem 2.5. First, when $\|r_j\|$ is large, ϵ_j is typically so small that there is no practical difference with requiring the preconditioning to be exact. Second, we do not know the residuals r_j in variable preconditioning—we have only the perturbed ones \underline{r}_j . Third, the corollary is based on the condition that r_j is monotonically decreasing, a too restrictive requirement. Hence we consider Corollary 2.6 to be mainly of theoretical interest.

3. Numerical examples. To empirically study when flexible preconditioning for BiCGStab is useful and to compare its performance with that of the counterparts of GMRES, we extend the example in Section 2 and perform more comprehensive experimentation. Consider the linear system arising from a discretization of the PDE (adopted from [23]):

$$-\Delta u + \gamma \mathbf{x} \cdot \nabla u + \beta u = f \quad (3.1)$$

with zero Dirichlet boundary condition. The linear system can be made indefinite and/or unsymmetric by changing the parameters γ and β . We discretized the 3D domain into a regular grid of size $n_1 \times n_2 \times n_3$ with spacing $h = 1/n_1$ and set $\gamma = 4/h$ to make the problem unsymmetric. We varied the parameter β to include both definite and indefinite cases. The right-hand side was chosen to be the vector of all ones, with the initial x_0 being zero. As a common practice in parallel solvers, if no specific preconditioner is mentioned, an iterative method (including the case of being used for inner iterations) was always preconditioned by block Jacobi/ILU(0), where each block was handled by one processor. The experiments were conducted on the supercomputer Titan (a Cray XK7 system) hosted at the Oak Ridge Leadership Computing Facility [22].

We tested several β values; Table 3.1 shows the results of three representative cases. As β decreases, the problem becomes harder to solve. In the first case ($\beta = 0.01/h^2$) the system is positive definite, but in the latter two cases ($\beta = -0.4/h^2$ and $-0.6/h^2$) the system is indefinite. Figure 3.1 shows the convergence history of BiCGStab and GMRES(30). These results were obtained on a $256 \times 256 \times 256$

grid using 16,384 processor cores. The relative residual tolerance and the maximum number of iterations were $1\text{e-}8$ and 200, respectively.

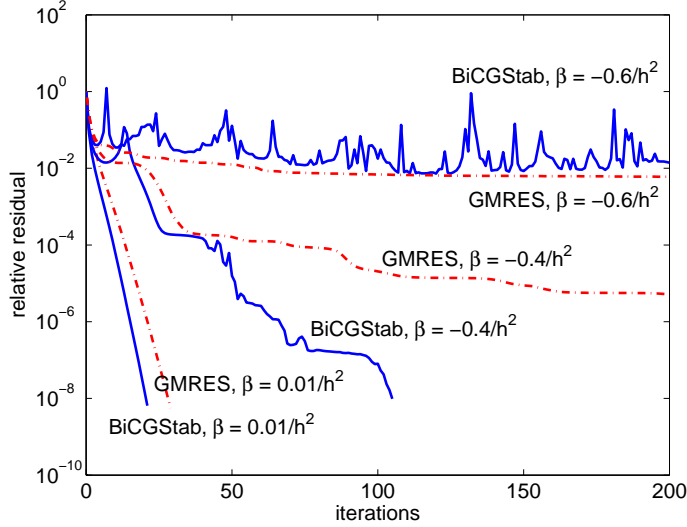


FIG. 3.1. *Convergence of BiCGStab and GMRES for three β 's (no inner iterations).*

We compared BiCGStab, FBiCGStab/BiCGStab, GMRES, and FGMRES/GMRES, where FBiCGStab/BiCGStab means FBiCGStab is preconditioned by BiCGStab, and similarly for FGMRES/GMRES. In both of the flexible methods, the stopping criterion for the inner iterations was either a residual tolerance or a fixed number of iterations. The two major comparison metrics are the number of matrix-vector multiplications (MatMult, column 4) and the number of calls to `MPI_Allreduce` (Allreduce, column 5). The number of floating point operations (Flops, column 6) for computing the inner products is listed for reference, but it is not a fair criterion for comparison because (i) the synchronization cost in `MPI_Allreduce` is much higher than that of floating point operations, and (ii) GMRES requires long-term recurrence (orthogonalization) whereas BiCGStab does not. The wallclock time is also listed, but note that the fluctuations in a distributed memory computing environment and the communication latency are factors that affect the actual running time.

Several observations are made according to Table 3.1 and other experiments with intermediate β values that are not shown in the table. First, for this set of test cases, the GMRES family in general performs better when the system is relatively easy to solve; but as the difficulty level increases, the BiCGStab family outperforms the GMRES family. One sees that in the first case ($\beta = 0.01/h^2$) GMRES is the best solver, whereas in the second case it failed to converge to the required tolerance but BiCGStab did converge. In the second case ($\beta = -0.4/h^2$), however, the fastest solver is still in the GMRES family. Nevertheless, when moving to the third case ($\beta = -0.6/h^2$), FBiCGStab with an inner tolerance stopping criterion clearly wins over all other solvers. This interesting phenomenon indicates that even for the same problem, different solvers may have significantly different performance profiles as the problem parameters vary.

Second, for a flexible Krylov method, using a fixed number of inner iterations can sometimes achieve excellent performance; but as the problem becomes harder

TABLE 3.1
Solution summary of (3.1) with three choices of β . Arrows point to the fastest run.

$\beta = 0.01/h^2$	Outer Iter.	Time (sec $\times 10^{-1}$)	MatMult	Inner Product	
				Allreduce	Flops $\times 10^9$
BiCGStab	21	1.132	42	84	3.010
FBiCGStab, inner rtol = 1e-3	2	1.397	58	120	4.246
FBiCGStab, inner rtol = 1e-2	2	1.335	44	92	3.233
FBiCGStab, inner rtol = 1e-1	4	1.364	54	116	4.037
FBiCGStab, inner max_it = 5	2	1.131	44	92	3.233
FBiCGStab, inner max_it = 10	1	1.147	40	82	2.910
FBiCGStab, inner max_it = 20	1	1.226	46	94	3.334
GMRES	29	0.874	29	58	9.494
FGMRES, inner rtol = 1e-3	3	1.122	40	83	6.081
FGMRES, inner rtol = 1e-2	4	7.200	37	78	4.240
FGMRES, inner rtol = 1e-1	7	1.153	38	83	3.292
FGMRES, inner max_it = 5	6	1.072	36	78	3.130
FGMRES, inner max_it = 10	3	1.066	33	69	4.237
FGMRES, inner max_it = 20	2	1.065	40	82	8.720
<hr/>					
$\beta = -0.4/h^2$		$\times 10^0$			$\times 10^6$
BiCGStab	105	0.356	210	420	1.505
FBiCGStab, inner rtol = 1e-3	2	0.668	438	880	3.144
FBiCGStab, inner rtol = 1e-2	2	0.499	270	544	1.948
FBiCGStab, inner rtol = 1e-1	6	0.832	586	1184	4.225
FBiCGStab, inner max_it = 5	fail	-	-	-	-
FBiCGStab, inner max_it = 10	fail	-	-	-	-
FBiCGStab, inner max_it = 20	25	2.517	2050	4150	14.800
GMRES	fail	-	-	-	-
FGMRES, inner rtol = 1e-3	3	1.556	1165	2298	37.890
FGMRES, inner rtol = 1e-2	4	1.062	803	1586	25.550
FGMRES, inner rtol = 1e-1	7	0.791	560	1113	17.160
FGMRES, inner max_it = 5	26	0.282	156	338	1.893
FGMRES, inner max_it = 10	14	0.278	154	322	2.134
FGMRES, inner max_it = 20	12	0.373	252	516	5.861
<hr/>					
$\beta = -0.6/h^2$		$\times 10^1$			$\times 10^7$
BiCGStab	fail	-	-	-	-
FBiCGStab, inner rtol = 1e-3	2	0.951	8314	16632	5.962
FBiCGStab, inner rtol = 1e-2	2	0.572	5024	10052	3.605
FBiCGStab, inner rtol = 1e-1	15	3.564	30120	60270	21.560
FBiCGStab, inner max_it = 5	fail	-	-	-	-
FBiCGStab, inner max_it = 10	fail	-	-	-	-
FBiCGStab, inner max_it = 20	fail	-	-	-	-
GMRES	fail	-	-	-	-
FGMRES, inner rtol = 1e-3	5	5.291	51670	101680	169.600
FGMRES, inner rtol = 1e-2	5	4.589	41800	82259	136.560
FGMRES, inner rtol = 1e-1	7	3.903	37757	74305	123.730
FGMRES, inner max_it = 5	fail	-	-	-	-
FGMRES, inner max_it = 10	fail	-	-	-	-
FGMRES, inner max_it = 20	fail	-	-	-	-

and harder, it is difficult to specify an appropriate number *a priori* to ensure the convergence of the outer iterations. One sees that in the first case, using a fixed number of inner iterations as the stopping criterion is in general preferable over using a residual tolerance. This is also true in the second case for the FGMRES solvers; however, for the FBiCGStab solvers the situation is completely opposite. In the third case, none of the solvers using a fixed number of inner iterations converged. In this sense, setting an inner tolerance is a more robust practice.

Third, one can choose an “optimal” inner tolerance for a flexible method. One sees that for FBiCGStab the inner tolerance $1\text{e-}2$ yields the best results in all the cases, whereas for FGMRES the tolerance is $1\text{e-}1$. This is consistent with the observation made in the experiments of flexible QMR [33], which states that the total solver cost first decreases, then increases as the inner solves are more and more exact. The “optimal” inner tolerance may be related to the convergence behavior of the inner iterations. One sees that in Figure 3.1 the relative residual norm of BiCGStab has a steep decrease at the beginning, until between $1\text{e-}1$ and $1\text{e-}2$. This may be the stopping point when BiCGStab becomes the most effective as an inner solve.

Fourth, the inner-outer iterations (that is, FBiCGStab/BiCGStab and FGMRES/GMRES) are often a better alternative to the standard iterations (that is, BiCGStab and GMRES). In harder problems the standard iterations did not converge whereas the inner-outer iterations did. In fact, the outer iterations converge extremely fast when an appropriate inner tolerance is used.

4. Application. In this section, we explore the use of flexible BiCGStab to improve the run-time performance of a large-scale application, PFLOTRAN [17, 19], where historically standard BiCGStab with block Jacobi/ILU(0) preconditioner has been the preferred linear solver. Based on the analysis in Section 2, we compose a multigrid (MG) preconditioner. Its use with BiCGStab yields two to three times improvement in solution time on $O(10^4)$ – $O(10^5)$ processor cores. When the coarse grid solver varies slightly (thus becoming a variable preconditioner and having to be used with FBiCGStab), we gain an additional 10–20% in overall runtime improvement.

PFLOTRAN is a state-of-the-art code for simulating multiscale, multiphase, multicomponent flow and reactive transport in geologic media. It solves a coupled system of mass and energy conservation equations for a number of phases, including air, water, supercritical CO_2 and a number of chemical components. The code utilizes finite volume or mimetic finite difference spatial discretizations and backward-Euler (fully implicit) timestepping. At each time step, Newton-Krylov methods are used for solving the resulting nonlinear algebraic equations. PFLOTRAN is built on the PETSc library [3, 4] and makes extensive use of PETSc iterative nonlinear and linear solvers.

The governing equations are described by Richards’ equation:

$$\frac{\partial}{\partial t}(\varphi s \rho) + \nabla \cdot \rho \mathbf{u} = \mathcal{S},$$

where φ denotes the porosity of the geologic medium, s the saturation (fraction of pore volume filled with liquid water), ρ the fluid density, \mathcal{S} a source/sink term representing water injection/extraction, and \mathbf{u} the Darcy velocity defined as

$$\mathbf{u} = -\frac{\kappa \kappa_r}{\mu} \nabla (P - \rho g z),$$

where P denotes fluid pressure, μ viscosity, κ the absolute permeability of the medium,

κ_r the relative permeability of water to air, g the acceleration of gravity, and z the vertical distance from a datum.

We consider two benchmark problems [18]:

Case 1: Cubic domain with a central injection well. This case models a $100\text{m} \times 100\text{m} \times 100\text{m}$ domain with a uniform effective permeability of 1 darcy and an injection well at the exact center.

Case 2: Regional flow without well near river. This case models a $5000\text{m} \times 2500\text{m} \times 100\text{m}$ region with a river at the eastern boundary.

In our numerical experiments, we ran each test case with a minimum number of time steps—six for case 1 and two for case 2—at which the number of linear iterations per Newton step became stabilized.

BiCGStab preconditioned by block Jacobi/ILU(0) has been the preferred linear solver for PFLOTTRAN because of its small memory consumption, compared with that of GMRES, and the empirically fast convergence. Similar to other Krylov methods, however, as the size of the application and the number of processors increase, BiCGStab encounters a well-known scaling difficulty for over 10,000 processor cores because of the bottleneck in vector inner-product calculations. The synchronization cost in `MPI_Allreduce` for computing inner products can constitute more than half of the solution time. In order to overcome the scaling difficulty, two algorithms were explored to successfully reduce the synchronization cost [18]. One was the use of flexible GMRES with a hierarchical GMRES preconditioner, where inner GMRES iterations are run on diagonal submatrices over subgroups of processor cores only. On each submatrix, inner GMRES is further preconditioned by block Jacobi/ILU(0). We denote this approach FGMRES/h-GMRES. The second algorithm was to use the IBiCGStab algorithm [39] with a Chebyshev preconditioner. We denote this approach to be IBiCGStab/Chebyshev. IBiCGStab is mathematically equivalent to BiCGStab, but the iterates are reorganized so that several inner products are computed together to reduce the number of synchronizations. The Chebyshev iterations may not be as effective in reducing the condition number as other Krylov iterations, but an advantage is that it does not require any inner product calculations and thus can be used as a fixed preconditioner. The two algorithms are based on the idea of reducing expensive global inner products across the entire system by using cheaper local inner products (e.g., FGMRES/h-GMRES) or inner iterations that do not compute inner products at all (e.g., IBiCGStab/Chebyshev).

The success in [18] prompted us to explore similar ideas for BiCGStab. The first attempt was to use BiCGStab in place of GMRES in FGMRES/h-GMRES, such as FBiCGStab/h-GMRES and FBiCGStab/h-BiCGStab. The use of a small constant number of inner iterations generally failed for convergence; the failure is not surprising because a few iterations make the preconditioner vary too much, and FBiCGStab does not guarantee any form of outer convergence as does FGMRES (monotonic decrease of residuals regardless of the varying of preconditioner). On the other hand, when we used an inner tolerance as the stopping criterion, convergence is seen for a reasonably chosen value (see Figure 4.1). In fact, the general trend for the same tolerance is similar for whichever preconditioner is used, either h-GMRES or h-BiCGStab, and the smaller the tolerance, the closer the residual history is to a fixed preconditioner case. However, this convergence was obtained at the cost of a large number of inner iterations. The overall run time cannot compete with the simple strategy of using BiCGStab preconditioned by block Jacobi/ILU(0).

We therefore searched for a preconditioner that converges faster to the level of

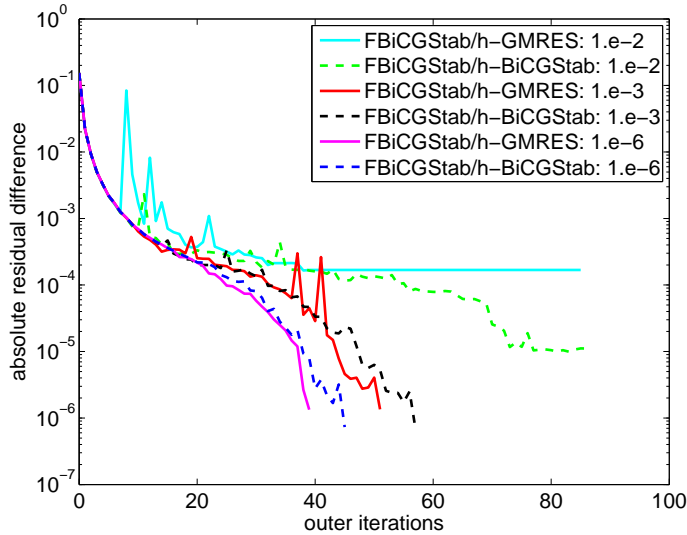


FIG. 4.1. *Convergence history of FBiCGStab/h-GMRES and FBiCGStab/h-BiCGStab for PFLOTRAN case 1 (mesh size $256 \times 256 \times 256$) using 512 processor cores.*

$1e-2$ to $1e-3$. Among all preconditioners examined, we settled on multigrid with a particular choice of the smoothers and the coarse-grid solver. Multigrid (MG) is generally applicable to steady-state or close to steady-state problems, such as the two cases of PFLOTRAN we consider here. In the standard setting, multigrid is a fixed preconditioner, because the smoothers (such as SOR or Chebyshev) require no inner-product calculations and the coarsest-grid problem is solved exactly by using a direct linear solver. In other words, one can, in principle, write down the linear operator for one cycle of multigrid and see that it does not vary.

For PFLOTRAN applications, the MG preconditioner was known to work well on a small number of processors, but the performance started tailing off at around 1,000 processor cores [16]. The difficulty to scale up the number of processors is that in the coarsest level, the problem is so small that the communication cost outweighs the computational cost. We cope with this difficulty by limiting the number of levels (effectively, three), such that the coarsest grid is not too small, and employing an iterative method to approximately solve the coarsest-grid problem. Two coarse grid solvers are (1) 100 Chebyshev iterations, and (2) 5 IBiCGStab iterations, each preconditioned by 20 Chebyshev iterations. For the former, even though the coarsest-grid problem is not solved to full accuracy, multigrid is still considered a fixed preconditioner because there are no inner product calculations. For the latter, clearly, multigrid is a variable preconditioner.

We conducted experiments with the above ideas on two computer systems: Intrepid, an IBM Blue Gene/P supercomputer located at the Argonne Leadership Computing Facility [1], and Titan, a Cray XK7 system located at the Oak Ridge Leadership Computing Facility [22]. Intrepid has 40,960 nodes, each consisting of one 850 MHz quad-core processor and 2GB RAM, resulting in a total of 163,840 cores, 80TB of memory, and a peak performance of 557 TFlops. Titan has 18,688 compute nodes, each consisting of one AMD 16-core Opteron 6274 processor running at 2.2GHz and 32GB of memory, giving a total of 299,008 cores.

TABLE 4.1
BiCGStab and FBiCGStab for PFLOTRAN on Intrepid (IBM Blue Gene/P), Case 1.

Number of Cores (Mesh Size)	IBiCGS/Cheby		BiCGStab/MG Smoother: Cheby CSolve: Cheby		FBiCGStab/Variable MG Smoother: Cheby CSolve: IBiCGS/Cheby	
	Iter.	Time	Iter.	Time	Iter.	Time (% Reduction)
512 (256x256x256)	547	212.8	29	97.2	23	86.4 (11%)
4,096 (512x512x512)	1006	365.1	43	121.3	33	106.1 (12%)
32,768 (1024x1024x1024)	1886	654.3	62	153.7	37	119.1 (23%)
163,840 (1600x1600x640)	2843	308.3	88	81.8	53	65.7 (20%)

TABLE 4.2
BiCGStab and FBiCGStab for PFLOTRAN on Intrepid (IBM Blue Gene/P), Case 2.

Number of Cores (Mesh Size)	IBiCGS/Cheby		BiCGStab/MG Smoother: Cheby CSolve: Cheby		FBiCGStab/Variable MG Smoother: Cheby CSolve: IBiCGS/Cheby	
	Iter.	Time	Iter.	Time	Iter.	Time (% Reduction)
16,384 (1600x816x320)	844	231.1	33	64.3	22	52.6 (18%)
98,304 (1600x1632x640)	1520	270.5	61	70.0	39	54.7 (22%)
163,840 (1600x1632x640)	1499	169.3	62	52.0	36	40.2 (23%)

Tables 4.1 through 4.4 compare the performance of (1) IBiCGStab/Chebyshev, (2) BiCGStab with the fixed MG preconditioner, and (3) FBiCGStab with the variable MG preconditioner on the two benchmark cases. The results of IBiCGStab/Chebyshev (columns 2–3) have been reported in [18] and are used here as the baseline of comparison. For MG preconditioners (columns 4–8), V-cycle was used; the smoothers were 2 steps of Chebyshev iterations because we found that Chebyshev outperformed other choices of smoothers. For both smoothers and coarse grid solvers, block Jacobi/ILU(0) was used as the innermost preconditioner.

Overall, FBiCGStab with an MG preconditioner is two to three times faster than the baseline IBiCGStab/Chebyshev on a large number of processor cores. Hence, the focus of comparison here is how much reduction in execution time one can achieve by using a variable preconditioner compared with using a fixed one. The percentage of reduction in execution time relative to the fixed MG preconditioner is listed in column 8 (% Reduction). As the size of the problem (hence coarsest grid) increases, the fixed number of iterations used for the coarsest-grid solver weakens the MG preconditioner, resulting in an increased number of outer iterations. Such an increase is less significant for the variable MG preconditioner because IBiCGStab/Chebyshev is more effective as a preconditioner than Chebyshev alone is. For example, when the number of cores becomes larger than 30,000, the number of outer iterations for using the variable MG preconditioner is almost half that of using the fixed MG preconditioner, leading to a reduction of 10% and 20% in overall execution time, on Titan and

TABLE 4.3
BiCGStab and FBiCGStab for PFLOTRAN on Titan (Cray XK7), Case 1.

Number of Cores (Mesh Size)	IBiCGS/Cheby		BiCGStab/MG Smoother: Cheby CSolve: Cheby		FBiCGStab/Variable MG Smoother: Cheby CSolve: IBiCGS/Cheby	
	Iter.	Time	Iter.	Time	Iter.	Time (% Reduction)
512 (256x256x256)	546	24.2	29	13.5	23	12.5 (7%)
4,096 (512x512x512)	1033	44.1	43	18.1	33	16.6 (8%)
32,768 (1024x1024x1024)	2073	89.0	62	27.4	37	24.3 (11%)
160,000 (1600x1600x640)	2407	52.0	91	24.9	55	22.5 (10%)

TABLE 4.4
BiCGStab and FBiCGStab for PFLOTRAN on Titan (Cray XK7), Case 2.

Number of Cores (Mesh Size)	IBiCGS/Cheby		BiCGStab/MG Smoother: Cheby CSolve: Cheby		FBiCGStab/Variable MG Smoother: Cheby CSolve: IBiCGS/Cheby	
	Iter.	Time	Iter.	Time	Iter.	Time (% Reduction)
1,600 (800x408x160)	411	18.0	20	8.8	20	8.8 (0%)
16,000 (1600x816x320)	829	29.6	30	11.3	22	10.9 (4%)
80,000 (1600x1632x640)	1578	53.1	60	16.7	38	15.0 (10%)
224,000 (1600x1632x640)	1501	20.9	66	16.6	37	14.8 (11%)

Intrepid, respectively.

Comparing the results obtained on the two machines, one sees a consistent iteration number (for some grid sizes, a slightly different number of processor cores was used across machines; this affected the innermost block Jacobi/ILU(0), thus making the iteration numbers slightly different). However, the time improvement of using a variable MG preconditioner is very different across the two machines. Because the clock rate of Intrepid is much lower than that of Titan, the solution time on Intrepid is longer. However, the communication network of Intrepid has a lower latency, so the global synchronization cost of `MPI.Allreduce` is smaller. We thus achieve better performance improvement on Intrepid because the variable MG preconditioner requires inner product calculations in the coarsest-grid solves. This phenomenon is not rare in practice. It showcases that for a solver, not only the theoretical convergence matters, but also the machine architecture plays an important role.

5. Concluding remarks. BiCGStab has been the de facto method of choice in many application domains for solving linear systems. Motivated by the challenges in large-scale scientific applications and extreme-scale computer architectures that encourage the use of variable preconditioners, we analyzed flexible BiCGStab and showed that the change of the convergence behavior with respect to standard BiCGStab is in accordance with the inaccuracy of the preconditioning solves. Thus, often a stop-

ping criterion with a moderate tolerance, say $1e-4$ to $1e-2$, for the preconditioning solves is favored in order both to maintain convergence and to reduce overall computation time. To this end, we demonstrated through numerical experiments, including the PFLOTRAN reacting flow application, that FBiCGStab with variable preconditioning yields superior performance on extreme-scale computers. This work provides insight on the practical use of FBiCGStab for large-scale applications.

Acknowledgments. We thank Satish Balay and Jed Brown for insightful discussions and assistance with experiments. The authors were supported by the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract DE-AC02-06CH11357.

REFERENCES

- [1] ALCF, *Intrepid Supercomputer*. <http://www.alcf.anl.gov/intrepid>.
- [2] J. ANG, K. EVANS, A. GEIST, M. HEROUX, P. HOVLAND, O. MARQUES, L. MCINNES, E. NG, AND S. WILD, *Report on the workshop on extreme-scale solvers: Transitions to future architectures*, Office of Advanced Scientific Computing Research, U.S. Department of Energy, 2012. Washington, DC, March 8-9, 2012.
- [3] S. BALAY, J. BROWN, K. BUSCHELMAN, V. EIJKHOUT, W. D. GROPP, D. KAUSHIK, M. G. KNEPLEY, L. C. MCINNES, B. F. SMITH, AND H. ZHANG, *PETSc users manual*, Tech. Rep. ANL-95/11 - Revision 3.3, Argonne National Laboratory, 2012.
- [4] S. BALAY, W. D. GROPP, L. C. MCINNES, AND B. F. SMITH, *Efficient management of parallelism in object oriented numerical software libraries*, in *Modern Software Tools in Scientific Computing*, E. Arge, A. M. Bruaset, and H. P. Langtangen, eds., Birkhauser Press, 1997, pp. 163–202.
- [5] P. G. BRIDGES, K. B. FERREIRA, M. A. HEROUX, AND M. HOEMMEN, *Fault-tolerant linear solvers via selective reliability*, CoRR, abs/1206.1390 (2012).
- [6] J. BROWN, M. G. KNEPLEY, D. A. MAY, L. C. MCINNES, AND B. F. SMITH, *Composable linear solvers for multiphysics*, in *Proceedings of the 11th International Symposium on Parallel and Distributed Computing (ISPD 2012)*, IEEE Computer Society, 2012, pp. 55–62. Also available as Argonne National Laboratory preprint ANL/MCS-P2017-0112.
- [7] A. CHRONOPOULOS AND C. W. GEAR, *S-step iterative methods for symmetric linear systems*, *Journal of Computational and Applied Mathematics*, 25 (1989), pp. 153–168.
- [8] H. A. V. DER VORST AND C. VUIK, *GMRESR: a family of nested GMRES methods*, *Numerical Linear Algebra with Applications*, 1 (1994), pp. 369–386.
- [9] A. EL MALIKI, R. GUENETTE, AND M. FORTIN, *An efficient hierarchical preconditioner for quadratic discretizations of finite element problems*, *Numerical Linear Algebra with Applications*, 18 (2011), pp. 789–803.
- [10] J. V. ESHOF AND G. L. G. SLEIJPEN, *Inexact Krylov subspace methods for linear systems*, *SIAM Journal on Matrix Analysis and Applications*, 26 (2004), pp. 125–153.
- [11] R. FLETCHER, *Conjugate gradient methods for indefinite systems*, *Lecture Notes in Mathematics*, 506 (1976), pp. 73–89.
- [12] P. GHYSELS, T. ASHBY, K. MEERBERGEN, AND W. VANROOSE, *Hiding global communication latency in the GMRES algorithm on massively parallel machines*, Tech. report 04.2012.1, Intel Exascale Lab, Leuven, Belgium, 2012.
- [13] E. GILADI, G. H. GOLUB, AND J. B. KELLER, *Inner and outer iterations for the Chebyshev algorithm*, *SIAM J. Numer. Anal.*, 35 (1995), pp. 300–319.
- [14] G. H. GOLUB AND Q. YE, *Inexact preconditioned conjugate gradient method with inner-outer iteration*, *SIAM Journal on Scientific Computing*, 21 (1999), pp. 1305–1320.
- [15] D. E. KEYES, L. C. MCINNES, C. WOODWARD, W. GROPP, E. MYRA, M. PERNICE, J. BELL, J. BROWN, A. CLO, J. CONNORS, E. CONSTANTINESCU, D. ESTEP, K. EVANS, C. FARHAT, A. HAKIM, G. HAMMOND, G. HANSEN, J. HILL, T. ISAAC, X. JIAO, K. JORDAN, D. KAUSHIK, E. KAXIRAS, A. KONIGES, K. LEE, A. LOTT, Q. LU, J. MAGERLEIN, R. MAXWELL, M. MCCOURT, M. MEHL, R. PAWLOWSKI, A. P. RANGLES, D. REYNOLDS, B. RIVIÈRE, U. RÜDE, T. SCHEIBE, J. SHADID, B. SHEEHAN, M. SHEPHARD, A. SIEGEL, B. SMITH, X. TANG, C. WILSON, AND B. WOHLMUTH, *Multiphysics Simulations: Challenges and Opportunities*, *International Journal of High Performance Computing Applications*, 27 (2013), pp. 4–83.

- [16] B. LEE AND G. HAMMOND, *Parallel performance of preconditioned Krylov solvers for the Richards equation*. Manuscript, 2010.
- [17] P. LIGHTNER ET AL., *PFLOTRAN project*. <http://ees.lanl.gov/pfлотran/>.
- [18] L. C. MCINNES, B. SMITH, H. ZHANG, AND R. T. MILLS, *Hierarchical and nested Krylov methods for extreme-scale computing*, Preprint ANL/MCS-P2097-0512, Argonne National Laboratory, 2012. submitted to Parallel Computing.
- [19] R. T. MILLS, V. SRIPATHI, G. MAHINTHAKUMAR, G. HAMMOND, P. C. LIGHTNER, AND B. F. SMITH, *Engineering PFLOTRAN for scalable performance on Cray XT and IBM BlueGene architectures*, in Proceedings of SciDAC 2010 Annual Meeting, 2010.
- [20] M. MOHIYUDDIN, M. HOEMMEN, J. DEMMEL, AND K. YELICK, *Minimizing communication in sparse matrix solvers*, in Proceedings of SC09, ACM, 2009.
- [21] Y. NOTAY, *Flexible conjugate gradients*, SIAM Journal on Scientific Computing, 22 (2000), pp. 1444–1460.
- [22] OLCF, *Jaguar Supercomputer*. <https://www.olcf.ornl.gov/computing-resources/jaguar/>.
- [23] Y. SAAD, *A flexible inner-outer preconditioned GMRES algorithm*, SIAM Journal on Scientific Computing, 14 (1993), pp. 461–469.
- [24] Y. SAAD, *Iterative Methods for Sparse Linear Systems, 2nd edition*, SIAM, Philadelphia, PA, 2003.
- [25] Y. SAAD AND M. H. SCHULTZ, *GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Stat. Comput., 7 (1986), pp. 856–869.
- [26] Y. SAAD AND M. SOSONKINA, *pARMS: A package for the parallel iterative solution of general large sparse linear systems user's guide*, Tech. Rep. UMSI2004-8, Minnesota Supercomputer Institute, University of Minnesota, 2004.
- [27] V. SIMONCINI AND D. SZYLD, *Flexible inner-outer Krylov subspace methods*, SIAM Journal on Numerical Analysis, (2003), pp. 2219–2239.
- [28] V. SIMONCINI AND D. B. SZYLD, *Theory of inexact Krylov subspace methods and applications to scientific computing*, SIAM Journal on Scientific Computing, 25 (2003), pp. 454–477.
- [29] G. L. SLEIJPEN, P. SONNEVELD, AND M. B. VAN GIJZEN, *Bi-CGSTAB as an induced dimension reduction method*, Applied Numerical Mathematics, 60 (2010), pp. 1100–1114.
- [30] G. L. SLEIJPEN AND M. B. VAN GIJZEN, *Exploiting BiCGstab(ℓ) strategies to induce dimension reduction*, SIAM J. Sci. Comput., 32 (2010), pp. 2687–2709.
- [31] P. SONNEVELD AND M. B. VAN GIJZEN, *IDR(s): A family of simple and fast algorithms for solving large nonsymmetric systems of linear equations*, SIAM Journal on Scientific Computing, 31 (2008), pp. 1035–1062.
- [32] E. D. STURLER AND H. A. VAN DER VORST, *Reducing the effect of global communication in GMRES(m) and CG on parallel distributed memory computers*, Applied Numerical Mathematics, 18 (1995), pp. 441–459.
- [33] D. B. SZYLD AND J. A. VOGEL, *FQMR: A flexible quasi-minimal residual method with inexact preconditioning*, SIAM Journal on Scientific Computing, 23 (2001), pp. 363–380.
- [34] H. VAN DER VORST, *BiCGSTAB: A fast and smoothly converging variant of BiCG for the solution of nonsymmetric linear systems*, SIAM J. Sci. Stat. Comput., 13 (1992), pp. 631–644.
- [35] M. B. VAN GIJZEN, G. L. SLEIJPEN, AND J.-P. M. ZEMKE, *Flexible and multi-shift induced dimension reduction algorithms for solving large sparse linear systems*, Tech. Rep. 11-06, Delft University of Technology, 2011.
- [36] J. VAN ROSENDALE, *Minimizing inner product data dependencies in conjugate gradient iteration*, in Proceedings of the IEEE International Conference on Parallel Processing, IEEE Computer Society, 1983.
- [37] J. A. VOGEL, *Flexible BiCG and flexible Bi-CGSTAB for nonsymmetric linear systems*, Applied Mathematics and Computation, 188 (2007), pp. 226–233.
- [38] R. VUDUC, *Quantitative performance modeling of scientific computations and creating locality in numerical algorithms*, PhD thesis, Massachusetts Institute of Technology, 1995.
- [39] L. T. YANG AND R. BRENT, *The improved BiCGStab method for large and sparse unsymmetric linear systems on parallel distributed memory architectures*, in Proceedings of the Fifth International Conference on Algorithms and Architectures for Parallel Processing, IEEE, 2002.

Government License. The submitted manuscript has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract

No. DE-AC02-06CH11357. The U.S. Government retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.